

METHOD FOR COMPILING PROGRAM COMPONENTS  
IN A MIXED STATIC AND DYNAMIC ENVIRONMENT

Background Of The Invention

5

Field of the Invention

10

The present invention relates to computer programming. More specifically, the invention relates to a method and variant of the method to compile programs or components of a program in a mixed static and dynamic environment.

Background Description

15

Most programming languages use the concept of a *data type* to identify a set of objects and operations that may be performed on those objects. Data types may be *primitive* (built into the language) or *user-defined*. A *class* in a programming language is used to create a user-defined type. A program written in an object-oriented manner can be viewed as a collection of classes. Classes contain declarations of both data and executable code in the form of *methods*. Herein, such methods are referred to as *procedures*.

25

30

Some programming languages, for example, Java®, have dynamic features like run-time binding of method calls and dynamic class loading. The term *virtual machine* is used herein to refer to the execution environment of such programming languages. Implementing a virtual machine for such a language may involve either interpreting the program or compiling it into the native code of the



target machine. Because interpretation incurs a high run-time overhead, virtual machines often rely on compilation for delivering high performance. Two prominent approaches to compilation in a virtual machine are dynamic compilation and static compilation. Dynamic compilation involves performing the translation of a program component (such as method or a collection of methods) to native machine code at run-time, before executing that program component. Static compilation involves performing the translation in an offline manner and generating one or more binary codes to be executed at run-time. Examples of virtual machines for Java using dynamic compilation include the IBM DK and the Sun JDK. Examples of static compilers for a Java-like language include JOVE, Tower Technologies TowerJ and the NaturalBridge BulletTrain compilers.

There are many problems with existing approaches to implementing virtual machines for dynamic languages. The problems with dynamic compilation include:

1. *Performance overhead of compilation at run-time:* The overhead of compilation is incurred every time the program is executed and is reflected in the overall execution time. Therefore, dynamic compilers tend to be less aggressive in applying optimizations that require deep analysis of the program.

2. *Testability and serviceability problems of the generated code:* Dynamic compilers that make use of run-time information about data characteristics to drive optimizations can lead to a different binary executable

being produced each time the program is executed. This can create reliability problems, as the code being executed may never have been tested.

5        3. *Large memory footprint:* A dynamic compiler is a complex software system with several interacting components, particularly if it supports aggressive optimizations. Hence, it usually has a large memory footprint, which gets directly added to the memory footprint of the application, since the dynamic compiler is invoked at run time. The memory footprint is particularly important for embedded systems, where the memory available on the device is limited.

10       Static compilation for dynamic languages leads to the following problems:

15       1. *Dynamic binding:* The code for dynamically linked class libraries may not be available during static compilation of a program, causing opportunities for interprocedural optimizations to be missed. Furthermore, the rules for binary compatibility in dynamic language like Java make it illegal to apply even simple inter-class optimizations -- e.g., method inlining across class boundaries -- unless the system has the ability to undo those optimizations in the event of changes to other classes.

20       2. *Dynamic class loading:* In general, dynamic class loading, as defined in languages like Java, requires the ability to handle a sequence of bytecodes representing a class (not seen earlier by the compiler) at run time. Hence, it is impossible for a virtual machine to support

a feature like dynamic class loading with a pure static compiler.

5 A digest of a data stream is a one-way hash function of the contents of the data stream that, with a very high probability, yields a different value if there are any changes made to the contents of the data stream. The Java 2 Security API supports secure hash functions to obtain the digest of a data stream or message.

10 Prior art for reducing the cost of dynamic compilation of Java can be found in An annotation-aware Java virtual machine implementation, *Proc. ACM SIGPLAN 1999 Java Grande Conference*, June 1999, A. Azevedo, A. Nicolau, and 15 J. Hummel. The AJIT compiler annotates the byte-code with machine independent analysis information that allows the JIT to perform some optimizations without having to dynamically perform analysis. A serious limitation of this system is that program transformation and code 20 generation still occur at application execution time.

25 Prior art for reducing the cost of dynamic compilation can be found in A general approach for run-time specialization and its application to C, *23rd ACM SIGACT-SIGPLAN Symposium on the Principles of Programming Languages*, pages 145-156, January 1996, C. Consel and F. Noel; An evaluation of staged run-time optimizations in DyC, *Proceedings of ACM SIGPLAN Conference on Programming Language Design and Implementation*, May 1999, B. Grant, 30 M. Philipose, M. Mock, C. Chambers, and S. Eggers; and Dynamic specialization in the Fabius system, *ACM Computing Surveys*, September 1998, M. Leone and P. Lee.

DyC is a selective dynamic compilation system for C, which reduces the dynamic compilation overhead by statically preplanning the dynamic optimizations. Based on user annotations that identify variables with relatively few run-time values, it applies partial evaluation techniques to partition computations in regions affected by those variables into static and dynamic computations.

Other systems, such as Tempo (see, A general approach for run-time specialization and its application to C, 23rd ACM SIGACT-SIGPLAN Symposium on the Principles of Programming Languages, pages 145-156, January 1996, C. Consel and F. Noel.) and Fabius (see, Dynamic specialization in the Fabius system, ACM Computing Surveys, September 1998, M. Leone and P. Lee.) support a similar staging of optimizations based on user annotations. All of these approaches have several limitations. First, they are unable to apply the staging of optimizations in the absence of user annotations. Second, they still require substantial code generation at run-time, and can only save the overhead of a few compiler optimizations. Third, they do not perform security checks to ensure the validity of code. Further, they do not deal with languages like Java, which have many dynamic features that make it difficult to generate code ahead of execution.

Prior art for recording the persistent execution state of a virtual machine for the Java platform can be found in Orthogonal persistence for the Java platform: Draft specification, October 1999,

Sub A' 7

<http://www.sun.com/research/forest/index.html>, M. Jordan and M. Atkinson; and Persistent execution state of a Java Virtual Machine, *Proc. ACM 2000 Java Grande Conference*, San Francisco, June 2000, T. Suezawa. These systems provide support for checkpointing the state of a Java application and virtual machine. They do not store the executable code for various procedures.

#### SUMMARY OF THE INVENTION

An object of this invention is to provide an improved method of compiling programs or components of a program in a mixed static and dynamic environment.

Another object of the present invention is to reduce the amount of time and memory spent in run-time compilation, while strictly honoring the semantics of dynamic features of a programming language.

A further object of this invention is to provide an improved method of compiling programs or components of a program in a mixed static and dynamic environment, so as to reduce the amount of time and resources spent in run-time compilations, and so as to exercise greater control over testing of the executable code for the program.

These and other objectives are attained with a method and system for a virtual machine in which compilation of a procedure is performed by (A) generating a persistent image, ahead of run time, that contains code for that procedure, and performing the following steps at run time; (B) checking for the existence and validity of a

code image for said procedure; (C) adapting the code  
image to the current execution context; and (D) using  
run-time compilation of the procedure if its code image  
does not exist, is invalid, or cannot be successfully  
adapted to the new execution context.

The preferred embodiment of this invention, as described  
below in detail, allows global interprocedural  
optimizations to be performed on the program, even if the  
programming language supports dynamic binding. Variants  
of the method show how one or several of the features of  
the method may be performed. The invention is  
particularly useful in the context of implementing Java  
Virtual Machines, although it can also be used in  
implementing other programming languages.

#### BRIEF DESCRIPTION OF THE DRAWINGS

The foregoing and other objects, aspects and advantages  
will be better understood from the following detailed  
description of a preferred embodiment of the invention  
with reference to the drawings, in which:

Figure 1 shows a block diagram of a prior art virtual  
machine in the context of which this invention may be  
used.

Figures 2 and 3 show a block diagram of a virtual machine  
using a method of this invention.

Figure 4 shows a block diagram of a QSI writer.  
Figure 5 shows a block diagram of the dependence  
recorder.

Figure 6 shows pseudocode for the adaptation annotation recorder component.

Figure 7 shows a flow chart of the QSI recorder.

Figure 8 shows a flowchart of the QSI repository system.

Figures 9 and 10 shows a flow-chart of the QSRT compiler component.

Figure 11 shows a flowchart describing the validation checks performed on a QSI.

Figure 12 shows pseudocode for a method of adapting the code and auxiliary information for a procedure to a new execution context.

## DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENTS OF THE INVENTION

### Using a Mixed Static and Dynamic Environment

Figure 1 shows a prior art system, a virtual machine, to which this invention is applied. A computer program (100) is transformed into executable code (102) by the compiler (101). The compiler may either be invoked at run-time or in an offline manner. The executable code is run by a run-time system (103).

Figures 2 and 3 show a system using an embodiment of this invention. In the preferred embodiment, the compilation activity is broken up into two phases, described in



Figures 2 and 3. Referring now to Figure 2, the computer program (200) is processed by a quasi-static image generator compiler (201), referred to as *QSI writer*. The *QSI writer* produces one or more quasi-static images (202), referred to as *QSI's*, which are persistent images of the executable code. The *QSI's* are stored for subsequent use by the virtual machine using a *QSI repository system* (203). Referring to Figure 3, the computer program (304), in the form of source code or intermediate language code, such as *bytecode* in a Java environment, is processed at run-time by a quasi-static run-time compiler (305), referred to as *QSRT compiler*. The *QSRT compiler* uses the *QSI repository system* (306) (which is identical to 203 in Figure 2) to retrieve the *QSI's* (307) containing executable codes for various components of the program. After processing the *QSI*, the *QSRT compiler* generates executable code (308) that is used by the run-time system (309) for executing the program.

#### Generation of Quasi-Static Images

Figure 4 shows a block diagram of a *QSI writer*. In the preferred embodiment of the method, the *QSI writer* is obtained by modifying a run-time compiler from prior art. In another embodiment, it is obtained by modifying a static, offline compiler. A front-end (401) processes the program to produce an intermediate code representation (402), which is fed to an optimizer (403) that produces optimized intermediate code (404). The front-end and optimizer represent well-known components of prior art compilers, and may be organized in different

Sub A<sup>2</sup> 7

ways, including, being organized in the form of multiple modules. The method of this invention adds to the optimizer a component (405) to record dependencies between different modules. This component is described further in Figure 5. The optimized intermediate code annotated with dependence information (406) produced by this component is processed by the back-end code generator (407), which may ignore the annotations on dependence information in the process of producing executable binary code (408). The method of this invention adds to the code generator an *adaptation annotation recorder* component (409), described further in Figure 6, which produces a further annotated executable code (410) with annotations to help adapt the code to a new execution context. The *QSI recorder* (411), described further in Figure 7 produces the QSI (412) which is stored for later processing.

The dependence recorder (405) from Figure 4 is described further in Figure 5. The *fine-grain dependence recorder* (500) keeps track of global optimizations performed by the compiler, and produces a list of fine-grain dependencies (501). In the preferred embodiment, these dependencies are recorded in the form of *class to procedure* dependencies. While compiling a procedure *A.foo* (i.e., a procedure *foo* of class *A*), for each optimization that exploits some information from a different class *B* (e.g., if a method *B.moo* is inlined into *A.foo*), the fine-grain dependence recorder in the method adds class *B* to the set of classes on which the code for *A.foo* is dependent. This allows the compiler, as explained later in the description of the QSRT

compiler, to perform dependence checks during program execution to avoid using stale code for a procedure in the event of changes to other code on which the code for that procedure is dependent.

5

The fine-grain dependencies (501) are processed by a *dependence granularity adjuster* (502) which replaces some fine-grain dependencies by coarser-grain dependencies to produce the final list of dependence annotations (503) to be used in the QSI. In the preferred embodiment, the dependence granularity adjuster examines the dependencies recorded for various procedures of each class. It factors out the dependencies that are common to all procedures in the class and records them at the *class to class* dependence level. The remaining dependencies continue to be recorded at the class to method level. For example, if a class A has two methods *foo* and *bar*, which are dependent, respectively, on classes B, C and B, D, the compiler would record the dependence of class A as a whole on B, and additionally, the dependence of *foo* on C and of *bar* on D.

10

15

20

25

30

The adaptation annotation recorder component (409) from Figure 4 is described further in the pseudocode shown in Figure 6. The list of annotations is initialized to be empty in Line 601. The loop in Line 602 processes each instruction in the machine code and each item recorded in the auxiliary information, such as exception tables and garbage collection maps in a language supporting exceptions and garbage collection. Line 603 checks if the current instruction or item is dependent on the



an instruction in a method *foo* of class *bar* which loads the static field *stats.count* (field *count* of class *stats*). The following PowerPC load instruction is generated to access the field:

5       lwz   R1=@{JTOC + offset of field *stats.count*  
JTOC is a dedicated register pointing to the table of global variables, and offset of field *stats.count* is an immediate-signed field giving the position of *stats.count* in that table. The value of the offset is assigned when  
10       the class *stats* is loaded. The adaptation annotation for the instruction is <I, T, S>, where I is the offset of the *lwz* instruction, T is an identifier denoting static field access, and S is the symbolic reference to the  
15       *constant pool entry* (as defined in the Java language specification (see *The Java Language Specification (Java Series)*, James Gosling, Bill Joy and Guy L. Steele, Jr. Addison-Wesley Publishing Company, Reading, MA.)) in the class *bar* for *stats.count*. Due to procedure inlining across class boundaries, a procedure may contain  
20       references that do not appear in the constant pool of its defining class. The compiler in the preferred embodiment creates an *extended constant pool* to handle relocation for references that are imported from other classes. An  
25       extended constant pool entry consists of the pair <N, C>, where N is an index into the constant pool of the class C from which this reference has been imported.

It should be noted that a static field reference is used only as an example to illustrate how adaptation  
30       annotations are recorded. There are other kinds of instructions that need annotations for the purpose of adaptation. For example, loads and stores of instance





In accordance with the layout of QSI shown in Figure 8, 705 leaves space for recording a digital signature for the QSI. The class to class dependence information, as obtained by the dependence recorder (405) described earlier, is written in 706, as a list of other classes on which the code being recorded in this QSI is dependent. A directory containing pointers to various procedure codes is created in 707. Note that a virtual machine may decide to create more than one code version for a given procedure, in order to perform optimizations based on specialization of procedures. The code for each procedure, along with auxiliary information such as exception tables, garbage collection maps, dependence information on other classes, and annotations for adaptation to a new execution context, is written to the QSI in 708.

A digest of the contents of the QSI is computed using a predetermined secure hashing function [see Proposed Federal Information Processing Standard for Secure Hash Standard. *Federal Register*, 57 (21), pages 3747-3749, January 1992). The digest is then encrypted using a well-known method (see *Applied Cryptography: Protocols, Algorithms, and Source Code in C*, B. Schneier, John Wiley and Sons, 1996.) to obtain a digital signature for the QSI, which is recorded at its predefined place in QSI in 709. The digital signature enables the virtual machine to detect any tampering of the QSI by a (malicious) user. Finally, 710 ensures that the above process (comprising steps 700 through 709) is continued for each class that is loaded during program execution.



09521571.072100

A flowchart of the QSI repository system (306) from Figure 3 is shown in Figure 8. A QSI may be stored in a file or in the memory itself. The preferred embodiment uses a file. The first step (800) is to identify where to place the QSI for a class. The QSI may be logically viewed as a part of the file containing the code for class seen by the virtual machine. In the preferred embodiment, the QSI is stored in a separate file with a .qsi suffix, but the method keeps track of the association between each QSI file and the original file containing the class code. The QSI machine uses a definite mapping (in step 800) to determine the directory in which a QSI is placed, given a unique identification of the class.

How this mapping is used may be illustrated with an example from a Java virtual machine. The location in which a classfile is stored in Java can be viewed as having two components: the *repository* containing the class, and the directory structure implied by the *fully qualified name* of the class [8]. For example, a class `MyPackage.Foo`, appearing in a repository `/vol/jdk/classes` on an AIX platform, is stored in the directory `/vol/jdk/classes/MyPackage`. The repository containing a class is identified by its defining class loader (e.g., using a search based on the classpath environment variable). For each class loader, a fixed mapping is defined from the name of the repository holding the class to the repository holding the QSI file, should it exist. Consider a class loader that loads classes over the network. The preferred method would use a local repository for the QSI files. Within a repository, the

method uses the same directory structure for a QSI file as that implied by the fully qualified name of the class. In the above example (for the class Foo in /vol/jdk/classes/MyPackage), given a QSI repository mapping function that replaces the string ``classes'' by the string ``qsi'', the corresponding QSI will be stored as Foo.qsi in the directory /vol/jdk/qsi/MyPackage.

Returning to Figure 8, step 801 checks for the existence of a QSI for the given class in the directory identified in step 800. A write request for a QSI is further processed using steps 802 through 803, while a read request for a QSI is processed using steps 804 through 805. Step 802 checks if the existing QSI should be modified in response to the write request. In the preferred embodiment, this checks the timestamp of the existing QSI. If it finds that the QSI is up-to-date, i.e., more recent than the file holding the class code, it decides not to overwrite the QSI. If the QSI is not up-to-date, then it deletes the older QSI, and the system proceeds to writing the new QSI in step 803. It should be noted that even when compiling a program for the first time, a QSI for a class from a library that is shared with other programs may already exist. Step 804 is followed for a read request if Step 801 shows that a QSI exists - it simply returns the QSI file. If for a read request, 801 shows that a QSI does not exist, a null value is returned by the read request in 805.

## Program Execution: Reuse of Quasi-Static Images

Figures 9 and 10 show a flow-chart of the QSRT compiler component (305) shown in Figure 3. The QSRT compiler is obtained by modifying a prior art run-time optimizing compiler in accordance with the methods of this invention. A controller (900) in the run-time compiler makes decisions on when a procedure is to be compiled. Step 901 checks if a QSI for the declaring class of that procedure has been already loaded. If it has not been loaded, 902 checks for the existence of a QSI for that class in the system, using a read request of the QSI repository system (306). If a non-null QSI is returned by the step, 903 performs validation and security checks on the QSI to determine if the QSI can safely be used. This step is described further in Figure 11. Step 904 reads the dependence list for the class stored in the QSI and performs dependence checks to see if any of the other classes, on which the code for the given class is dependent, have changed. This is done by comparing the timestamps of files holding the codes for classes in the dependence list with the timestamp of the given QSI. If any of those files have a more recent timestamp than the timestamp of the QSI, the dependence check fails.

Such dependence checking allows the virtual machine to ensure the validity of generated code while performing global optimizations like inlining across class boundaries, in the presence of changes to other codes. If the dependence check passes, 905 reads the method directory area in the QSI to look up the pointer to the code for the procedure to be compiled and reads that code

from the QSI. If this code is found, 906 performs dependence checks, again using timestamps, to see if any of the classes on which code for this method is dependent have changed since the time the method code was generated. If this check passes, 907 adapts the code and auxiliary information for the procedure to the new execution context. This step is further described in Figure 12. If the adaptation of code succeeds, 908 returns the adapted executable code as the compiled code for the procedure. If any of the previous steps 902 through 907 fail, as shown in the flowchart, 909 performs run-time compilation of the procedure.

Figure 11 shows a flowchart providing further details of step 903 in Figure 9. Step 1100 reads the header data from the QSI. Step 1101 performs a compatibility check to ensure that the recorded QSI has been produced under an environment compatible with the environment of the current virtual machine instance. This involves verifying the magic number recorded in the QSI, checking the virtual machine version, operating system version, and the target architecture identifier, and ensuring that those are compatible with the current virtual machine, operating system, and target architecture. Step 1102 performs an out-of-date check, using timestamps, to see if the QSI is older than the file holding code for the corresponding class. In a virtual machine for Java 2, this test includes an additional check to ensure that the defining class loader (if it is not the primordial class loader) for the class being compiled is identical to the defining class loader for the class at the time of QSI creation.





Publishing Company, Reading, MA.) in the class *bar* for *stats.count*. Note that the class *foo* is already loaded and resolved before the virtual machine controller decides to compile *bar*. Therefore, the virtual machine  
 5 has information mapping the entry *S* in the constant pool of class *foo* to a new entry in the table being used to hold static fields. The old offset for *stats.count* is replaced by the offset of this new entry.

10 Those skilled in the art will recognize that if *S* was referring to an entry in the *extended constant pool*, denoting an entry from the constant pool of another class whose method was inlined into *bar*, the adaptation could similarly be done using the mapping of constant pool for  
 15 that class. It should be further noted that in some situations, extra code has to be inserted to do the adaptation

In the above example, if the class *stats* has not been  
 20 loaded and resolved at the time of adaptation of the code for *foo.bar* (whereas the class was loaded and resolved when the QSI for *foo* was generated), the virtual machine would generate code of the following form:  
 label:

```

25      lwz    R1=@{JTOC + offset of field table (1)
      lwz    R1=@{R1 + field Id                      (2)
      if (R1 == 0) goto resolve                      (3)
      lwz    R1=@{JTOC + R1 // field access          (4)
      ...
  
```

```

30      resolve:
      (resolve field)                                (5)
  
```





In other alternate embodiments of the method, the generation of QSI files is not necessarily done in a separate phase from the execution of the program. In one such embodiment, the step 909 performing run-time compilation of a procedure is followed by addition of the newly generated code along with auxiliary information to an existing or new QSI for the class containing the procedure. This allows the QSI's for a code to evolve with time in response to adaptive compilation. However, this embodiment can lead to additional overhead at run-time to create the QSI.

An alternate embodiment uses digests of classfiles rather than using timestamps for out-of-date checks in step 1102. In this embodiment, the virtual machine records a digest of the original classfile and of the classfiles on the dependence list, in the QSI file. The out-of-date check is performed by comparing the digest of the current classfile with the recorded digest for that file. An advantage of this approach is that trivial changes to a classfile's timestamp (due to the file being *touched* or moved) do not cause an unnecessary invalidation of the quasi-static image file.

Yet another alternate embodiment performs the reading of procedure code from the QSI file in an *eager* manner rather than in a *lazy* manner in the QSRT compiler. Rather than reading the code and auxiliary information for a procedure from a QSI only when responding to a virtual machine request to compile a procedure (in Step 905), this information is read, in this embodiment, when

the QSI file is read for the first time, which in turn happens the first time that compilation is attempted of any method in that class. This is useful when most (or all) procedures stored in the QSI are used during program execution, because sequential I/O is more efficient due to buffer prefetching. On the other hand, this has a potential drawback of leading to unnecessary I/O if relatively few procedures stored in the QSI are needed during program execution.

An alternate embodiment allows multiple code versions of a procedure to appear in a single QSI. Therefore, library code may be specialized for different applications. Furthermore, since compilation is done in an offline manner, the compiler has more freedom to apply potentially expensive interprocedural analysis to discover opportunities for specialization.

An alternate embodiment uses a different strategy for recording dependence information in Step 405 to explore the trade-offs between the overhead of dependence checking and the likelihood of QSI invalidation during program execution time. In this embodiment, the compiler moves a dependence item shared by the important (but not all) procedures to a class-level dependence, to reduce the overhead of dependence checking. This comes at the expense of possibly invalidating the entire QSI file rather than invalidating just the quasi-static code for a single procedure. In yet another embodiment, the compiler keeps dependence information at a finer granularity (e.g., procedure-to-procedure dependence) to reduce the

chances of invalidating a QSI, at the expense of increased overhead of dependence checking.

Another embodiment of the method does not use run-time compilation in step 909, in order to handle procedures for which executable code could not be not obtained using a QSI. Rather than run-time compilation, the virtual machine uses interpretation of such a procedure. An advantage of this approach is that it leads to a smaller memory footprint at run-time, which can be particularly useful for embedded and hand-held devices.

Thus, it should be understood that the preferred embodiment is provided as an example and not as a limitation. While the invention has been described in terms of a single preferred embodiment, with several variants, those skilled in the art will recognize that the invention can be practiced with modification within the spirit and scope of the appended claims.